1 PhaeoEpiView: An epigenome browser of the newly assembled genome of the model diatom
2 *Phaeodactylum tricornutum*

3

4 Yue Wu[1][¥], Chaumier Timothée[1][¥], Eric Manirakiza[1], Alaguraj Veluchamy[2] and Leila Tirichine[1][*]

5 [1]Nantes Université, CNRS, US2B, UMR 6286, F-44000 Nantes, France
6
7 [2]St Jude Children's Research Hospital, Memphis, TN, US
8
9

10 [¥]Equal contribution

11

12

13 [*]Correspondence: tirichine-l@univ-nantes.fr; Tel.: +33-276645058

14

## Abstract

### Motivation

Recent advances in DNA sequencing technologies in particular of long reads type greatly improved genomes assembly leading to discrepancies between both published annotations and epigenome tracks which did not keep pace with new assemblies. This comprises the availability of accurate resources which penalizes the progress in research.

### Results

Here, we used the latest improved telomere to telomere assembly of the model pennate diatom *Phaeodactylum tricornutum* to lift over the gene models from Phatr3, a previously annotated reference genome. We used the lifted genome annotation including genes and transposable elements to map the epigenome landscape, namely DNA methylation and post translational modifications of histones providing the community with PhaeoEpiView, a browser that allows the visualization of epigenome data as well as transcripts on an updated reference genome to better understand the biological significance of the mapped data on contiguous genome rather than a fragmented one. We updated previously published histone marks with a more accurate mapping using monoclonal antibodies instead of polyclonal and deeper sequencing. PhaeoEpiView will be continuously updated with the newly published epigenomic data making it the largest and richest epigenome browser of any stramenopile. We expect that PhaeoEpiView will be a standard tool for the coming era of molecular environmental studies where epigenetics holds a place of choice.

### Availability

PhaeoEpiView is available at: https://PhaeoEpiView.univ-nantes.fr

**Introduction**

The genome of the model diatom *Phaeodactylum tricornutum* CCAP 1055/1 and the corresponding annotation were published in 2008 using whole genome shotgun paired-end Sanger sequencing (NCBI assembly ASM15095v2) (Bowler et al., 2008). Subsequently, Phatr3 annotation updated gene repertoire to introduce more than thousand novel genes and performed a comprehensive de novo annotation of repetitive elements showing novel classes of transposable elements using 90 RNA-Seq datasets combined with published expressed sequence tags and protein sequences (Rastogi et al., 2018). The first assembly of the genome contained 33 scaffolds among which 12 telomere-to-telomere chromosomes. Using long read sequencing, Filloramo et al., re-examined *P. tricornutum* assembly which led to additional sequence information but did not improve the continuity and chromosome-level scaffolds compared to the original reference genome (Filloramo et al., 2021). Recently, an approach combining long reads from the Oxford Nanopore minION platform and short high accurate reads from the Illumina NextSeq platform was used to perform a new assembly of *P. tricornutum* genome which led to 25 nuclear chromosomes improving thus the assembly (Giguere, 2021). However, Phatr3 annotation of *P. tricornutum* was not revisited in light of the new 25 telomeric chromosomes assembly which is often observed for most species where the annotations do not keep pace with new/improved assemblies.

*P. tricornutum* is an established model diatom widely used by an increasing community for fundamental research and biotech applications. Diatoms are one of the most abundant and highly diverse mostly photosynthetic eukaryotes, contributing to 20-25% of the Earth's global carbon dioxide fixation (Field et al., 1998) and their photosynthetic activity accounts for about 40% of the marine primary production (Armbrust, 2009; Falkowski et al., 1998). Diatoms are highly successful and widely spread occupying large territories including marine, freshwater, sea ice, snow and even moist terrestrial habitats.

While whole-genome sequencing is critical to better understanding the ecological success of diatoms, primary sequence is only the basis for understanding how to read genetic programs, another layer of heritable information superimposed on the DNA sequence is epigenetic information. It has already been proposed that the ecological success of phytoplankton is also due to the adaptive dynamics conferred by epigenetic regulation mechanisms because point mutation-based processes may be too slow to permit adaptation to a dynamic ocean

70   environment (Tirichine and Bowler, 2011). The epigenetic changes may lead to chromatin

71   modifications, which may cause a stable alteration in transcriptional activity even after

72   withdrawal of the triggering stress (Avramova, 2015). Pioneering work drew a comprehensive

73   map of epigenetic marks including several permissive and repressive PTMs and DNA

74   methylation in *P. tricornutum* and showed their contribution to mediate the response of diatom

75   cells to environmental factors (Rastogi et al., 2015; Veluchamy et al., 2013; Veluchamy et al.,

76   2015).

77   An important molecular tool box is available in *P. tricornutum* including epigenomic data

78   which are only found in the partial assembly. To make such a resource available on the newly

79   assembled genome, we used the new 25 to 25 telomere assembly to map the epigenetic data

80   including Post-translational modifications of histones (PTMs) and DNA methylation that were

81   previously published (Hoguin, 2021; Veluchamy *et al.*, 2013; Veluchamy *et al.*, 2015). Prior to

82   this, we lifted the Phatr3 annotation using a gene based approach.

83   **Implementation**

84   PhaeoEpiView was built using two steps (i) Phatr3 gene annotation lifting onto the new

85   25 chromosomes assembly (Phatr4) (ii) mapping of the previously published epigenetic marks

86   and transcripts on the new assembly. For more accuracy and homogeneity of the used data,

87   chromatin immunoprecipitation with deep sequencing was carried out using monoclonal

88   antibodies to replace two marks, H3K9me3 and H3K27me3 for which polyclonal antibodies

89   were used in the previous study.

90

91   In the first step, instead of whole genome-based comparison, we adapted a gene-based

92   sequence alignment for lifting the annotation from Phatr3 to Phatr4 assembly. Features such as

93   mRNA, CDS and exons from the reference Phatr3 were used to infer genes and transcripts in

94   target assembly. Using minimap2 and Liftoff tools, exons are aligned first to preserve the gene

95   structure of the Phatr3 annotation (Li, 2018; Shumate, 2021). Minimap is used with 50

96   secondary mappings, end bonus of 5 and chaining score of 0.5. Genes are lifted and considered

97   mapped successfully if the alignment coverage and sequence identity in the child features

98   (usually exons/CDS) is >= 50%. Unplaced genes and genes with extra copy number are tagged

99   and separated (Supplementary Table 1). Out of the 12178 genes from Phatr3 annotation, 11739

100  were lifted successfully (Supplementary File 2).

101

102    In order to validate the lifted annotation, we aligned RNAseq reads to both the previous
103    and the new genome version then compared every gene quantification. The vast majority of
104    them had a difference of quantification (Supplementary Figure 1) and length lower than
105    +/- 0.1% between Phatr3 and Phatr3_lift (Phatr4). Missing genes were then examined: 178 out
106    of 439 (40%) were found to be located on short regions that are no longer present in Phatr4
107    assembly according to whole-genome alignment provided in (Giguere, 2021), half of them
108    clustered on previous chr_5 and chr_21 (Supplementary Table 1). Most of the remaining 261
109    missing genes showed similarity to already lifted genes suggesting that they are either
110    duplicated or allelic.
111

112    In the second step, transposable elements annotation available from (Giguere, 2021) was
113    added to PhaeoEpiView as Phatr4 TEs track. Finally, previously published expression data at
114    two different time points, DNA methylation and PTMs tracks were implemented in the browser
115    and systematic comparison was made with the previous assembly mapping using unchanged
116    regions as anchors (Supplementary File 1). PhaeoEpiView was implemented as a Jbrowse2
117    instance (Buels et al., 2016) and made public on a virtual machine hosted at Nantes University
118    datacenter. It can currently display one track for each of the genes, TEs, transcript levels,
119    McrBC and Bisulfite-seq DNA methylation and five histone PTMs (H3K9/14Ac, H3K4me2,
120    H3K4me3, H3K9me2, H3K9me3, H3K27me3). The browser will be regularly updated with
121    relevant epigenomic data when published in the future, making PhaeoEpiView a live platform
122    for a comprehensive genomic and epigenomic resource of the model microalgae *P. tricornutum*.
123

## Conclusion

125    PhaeoEpiView is an open source browser that provides an up to date genome and
126    epigenome view of the model diatom *Phaeodactylum tricornutum*. With the lifted genes
127    annotation, the epigenome and transcriptome landscapes can be visualized on a fully assembled
128    genome providing an accurate view of epigenetic regulation of genes and TEs which was
129    incomplete on the previously fragmented genome. PhaeoEpiView allows users to upload their
130    own tracks in private session for visualization and data interpretation purposes. PhaeoEpiView
131    is intuitive, easy to use and represent the first epigenome browser of a photosynthetic unicellular
132    species which will undoubtedly contribute to boost research in microalgae and single celled
133    species in general.

## Acknowledgements

## References

Armbrust, E.V. (2009). The life of diatoms in the world's oceans. Nature *459*, 185-192. 10.1038/nature08057.

Bowler, C., Allen, A.E., Badger, J.H., Grimwood, J., Jabbari, K., Kuo, A., Maheswari, U., Martens, C., Maumus, F., Otillar, R.P., et al. (2008). The Phaeodactylum genome reveals the evolutionary history of diatom genomes. Nature *456*, 239-244. nature07410 [pii]

10.1038/nature07410.

Falkowski, P.G., Barber, R.T., and Smetacek, V.V. (1998). Biogeochemical Controls and Feedbacks on Ocean Primary Production. Science *281*, 200-207.

Field, C.B., Behrenfeld, M.J., Randerson, J.T., and Falkowski, P. (1998). Primary production of the biosphere: integrating terrestrial and oceanic components. Science *281*, 237-240.

Filloramo, G.V., Curtis, B.A., Blanche, E., and Archibald, J.M. (2021). Re-examination of two diatom reference genomes using long-read sequencing. BMC Genomics *22*, 379. 10.1186/s12864-021-07666-3.

Giguere, D.J., Bahcheli, A.T., Slattery, S.S., Patel, R.R., Flatley, M., Karas, B.J., Edgell, D.R., Gloo, G.B. (2021). Telomere-to-telomere genome assembly of Phaeodactylum tricornutum. doi: https://doi.org/10.1101/2021.05.04.442596.

Huang, R., Ding, J., Gao, K., Cruz de Carvalho, M.H., Tirichine, L., Bowler, C., and Lin, X. (2018). A Potential Role for Epigenetic Processes in the Acclimation Response to Elevated pCO2 in the Model Diatom Phaeodactylum tricornutum. Front Microbiol *9*, 3342. 10.3389/fmicb.2018.03342.

Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics *34*, 3094-3100. 10.1093/bioinformatics/bty191.

Rastogi, A., Maheswari, U., Dorrell, R.G., Vieira, F.R.J., Maumus, F., Kustka, A., McCarthy, J., Allen, A.E., Kersey, P., Bowler, C., and Tirichine, L. (2018). Integrative analysis of large scale transcriptome data draws a comprehensive landscape of Phaeodactylum tricornutum genome and evolutionary origin of diatoms. Sci Rep *8*, 4834. 10.1038/s41598-018-23106-x.

Shumate, A.a.S., S.L. (2021). Liftoff: Accurate Mapping of Gene Annotations. Bioinformatics *37*, 1639–1643.

Veluchamy, A., Lin, X., Maumus, F., Rivarola, M., Bhavsar, J., Creasy, T., O'Brien, K., Sengamalay, N.A., Tallon, L.J., Smith, A.D., et al. (2013). Insights into the role of DNA methylation in diatoms by genome-wide profiling in Phaeodactylum tricornutum. Nat Commun *4*. 10.1038/ncomms3091.

Veluchamy, A., Rastogi, A., Lin, X., Lombard, B., Murik, O., Thomas, Y., Dingli, F., Rivarola, M., Ott, S., Liu, X., et al. (2015). An integrative analysis of post-translational histone modifications in the marine diatom Phaeodactylum tricornutum. Genome Biol *16*, 102. 10.1186/s13059-015-0671-8.

177  Buels, R., Yao, E., Diesh, C.M., Hayes, R.D., Munoz-Torres, M., Helt, G., Goodstein, D.M.,
178  Elsik, C.G., Lewis, S.E., Stein, L., and Holmes, I.H. (2016). JBrowse: a dynamic web platform
179  for genome visualization and analysis. Genome Biol *17*, 66. 10.1186/s13059-016-0924-1.
180  Giguere, D.J., Bahcheli, A.T., Slattery, S.S., Patel, R.R., Flatley, M., Karas, B.J., Edgell, D.R.,
181  Gloo, G.B. (2021). Telomere-to-telomere genome assembly of Phaeodactylum tricornutum.
182  doi: https://doi.org/10.1101/2021.05.04.442596.
183  Hoguin, A., Ait Mohamed, O., Bowler, C., Genovesio, A., Vieira, F.R.J., Tirichine, L. (2021).
184  Evolutionary analysis of DNA methyltransferases in microeukaryotes: Insights from the model
185  diatom Phaeodactylum tricornutum. bioRxiv, doi: https://doi.org/10.1101/2021.06.11.447926.
186  Rastogi, A., Lin, X., Lombard, B., Loew, D., and Tirichine, L. (2015). Probing the evolutionary
187  history of epigenetic mechanisms: What can we learn from marine diatoms. AIMS Genetics *2*,
188  173-191. 10.3934/genet.2015.3.173.
189  Tirichine, L., and Bowler, C. (2011). Decoding algal genomes: tracing back the history of
190  photosynthetic life on Earth. Plant J *66*, 45-57. 10.1111/j.1365-313X.2011.04540.x.
191  Veluchamy, A., Lin, X., Maumus, F., Rivarola, M., Bhavsar, J., Creasy, T., O'Brien, K.,
192  Sengamalay, N.A., Tallon, L.J., Smith, A.D., et al. (2013). Insights into the role of DNA
193  methylation in diatoms by genome-wide profiling in Phaeodactylum tricornutum. Nat Commun
194  *4*. 10.1038/ncomms3091.
195  Veluchamy, A., Rastogi, A., Lin, X., Lombard, B., Murik, O., Thomas, Y., Dingli, F., Rivarola,
196  M., Ott, S., Liu, X., et al. (2015). An integrative analysis of post-translational histone
197  modifications in the marine diatom Phaeodactylum tricornutum. Genome Biol *16*, 102.
198  10.1186/s13059-015-0671-8.

199

200  **Supplementary methods**

201

202  **Culture and growth conditions**

203  *Phaeodactylum tricornutum* Bohlin Clone Pt1 8.6 (CCMP2561) cells were obtained from the

204  culture collection of the Provasoli-Guillard National Center for Culture of Marine

205  Phytoplankton (Bigelow Laboratory for Ocean Sciences, USA). Constantly shaken (100 rpm)

206  cultures were grown at 19˚C, 60 µmol photons $m^{-2}$ $s^{-1}$ and with a 12h light / 12h dark

207  photoperiod in sterile Enhanced Artificial Sea Water (EASW) medium (Vartanian, et al., 2009).

208  For Chromatin immunoprecipitation-sequencing, cultures were seeded at 50.000 cells/ml in

209  duplicate and grown side by side in 1000 ml erlens until early-exponential at $10^6$ cells/ml.

210  Culture growth was measured using a hematocytometer (Fisher Scientific, Pittsburgh, PA,

211  USA).

212

213  **Chromatin extraction and immunoprecipitation**

214  Chromatin isolation was performed as described previously (Lin, et al., 2012) with few

215  modifications. Briefly, the incubation step in buffer II is repeated several times until the pellet

216  becomes white. Each ChIP-seq experiment was conducted in two independent biological

217  replicates. Monoclonal antibodies from Cell Signaling Technology were used for

218    immunoprecipitation, H3K9me3 (13969), and H3K27me3 (9733).

219

220    **ChIP-Seq analysis**

221    Pair-end sequencing of H3K9me3 and H3K27me3 ChIP and input samples was performed on

222    Illumina NovaSeq platform with read length of 2 x 150 bp. Previously published ChIP

223    sequencing for H3K9me2, H3K9me3, H3K4me2, H3K27me3, H3K9/K14Ac and H3K4me3

224    were retrieved from NCBI's Gene Expression Omnibus accessions GSE68513 and GSE139676

225    (Veluchamy, et al., 2015; Zhao, et al., 2021). Raw reads were filtered and low-quality read pairs

226    were discarded using Trim Galore 0.6.7 (https://doi.org/10.5281/zenodo.5127899) with a read

227    quality (Phred score) cutoff of 20 and a stringency value of 3bp. Using the 25 to 25 telomere

228    assembly published in 2021 as a reference genome, the filtered reads were mapped using

229    Bowtie2 2.4.5 (Langmead and Salzberg, 2012). We then performed the processing and filtering

230    of the alignments using Samtools 1.15 "fixmate -m" and "markdup -r" modules (Danecek, et

231    al., 2021). Two biological replicates for each ChIP were performed and read counts showed a

232    good Pearson correlation by Deeptools multiBamSummary v3.5.1 with a bin size of 1000bp

233    (Ramirez, et al., 2014). To identify regions that were significantly enriched, we used MACS2

234    v2.2.7.1 (Zhang, et al., 2008) on the combination of the two replicates with "callpeak --qvalue

235    0.05 --nomodel --SPMR --bdg" options. In addition, extension size was set to the arithmetic

236    mean of the two IP replicates fragment size for each mark, as determined by MACS2 predicted

237    module with "-m 2 70" MFOLD value. Furthermore, "--broad" calling mode was activated for

238    H3K9me2 and H3K9me3 that were previously described as broad histone marks. For the

239    narrow marks H3K4me2, H3K9/K14Ac and H3K4me3, peaks summits were called with "--

240    call-summits". Following previously published work, SICER2 v1.0.3 (Zang, et al., 2009) was

241    used with "-w 200 -g 600 -fdr 0.05" to call peaks for H3K27me3.

242    Output normalized Fold Enrichment signal files were generated with MACS2 "bdgcmp"

243    module and transformed to BigWig using Deeptools bedGraphToBigWig. Then, Pearson

244    correlation between our new data and previously published data for H3K9me3 and H3K27me3

245    was performed using Deeptools plotCorrelation.

246

247    **Expression analysis**

248    Early and late exponential growth phase Illumina RNA-seq data (SRR5274697, SRR5274696,

249    SRR5274695 and SRR5274694) from (Murik, et al., 2019) were trimmed using Trim Galore

250    0.6.7 with a read quality (Phred score) cutoff of 20 and a stringency value of 3bp. Technical

251    replicates were merged and mapped to the reference assembly with STAR 2.7.10a

252 (https://www.ncbi.nlm.nih.gov/pubmed/23104886). Primary alignments only were processed

253 with Deeptools bamCoverage 3.5.0 with "--normalizeUsing BPM --ignoreDuplicates --

254 centerReads" to generate normalized coverage files to be displayed in PhaeoEpiView.

255

256 **DNA methylation analysis**

257 McrBC DNA methylation annotation data from

258 (https://www.nature.com/articles/ncomms3091) was lifted from the previous assembly to the

259 new 25 chromosomes using Liftoff. Bisulfite sequencing data (Hoguin, 2021) were processed

260 with Bismark v0.22.3 (https://pubmed.ncbi.nlm.nih.gov/21493656/) and methylated regions

261 having less than 50% methylated reads or less than 5 supporting reads were filtered out.

262

263

264 **References**

265

266 Danecek, P*., et al.* Twelve years of SAMtools and BCFtools. *Gigascience* 2021;10(2).

267 Hoguin, A., Ait Mohamed, O., Bowler, C., Genovesio, A., Vieira, F.R.J., Tirichine, L.

268 Evolutionary analysis of DNA methyltransferases in microeukaryotes: Insights from the model

269 diatom Phaeodactylum tricornutum. *bioRxiv, doi: https://doi.org/10.1101/2021.06.11.447926*

270 2021.

271 Langmead, B. and Salzberg, S.L. Fast gapped-read alignment with Bowtie 2. *Nat Methods*

272 2012;9(4):357-359.

273 Lin, X., Tirichine, L. and Bowler, C. Protocol: Chromatin immunoprecipitation (ChIP)

274 methodology to investigate histone modifications in two model diatom species. *Plant methods*

275 2012;8(1):48.

276 Murik, O*., et al.* Downregulation of mitochondrial alternative oxidase affects chloroplast

277 function, redox status and stress response in a marine diatom. *New Phytol* 2019;221(3):1303-

278 1316.

279 Ramirez, F*., et al.* deepTools: a flexible platform for exploring deep-sequencing data. *Nucleic*

280 *Acids Res* 2014;42(Web Server issue):W187-191.

281 Vartanian, M*., et al.* Plasticity and robustness of pattern formation in the model diatom

282 Phaeodactylum tricornutum. *The New phytologist* 2009;182(2):429-442.

283 Veluchamy, A*., et al.* An integrative analysis of post-translational histone modifications in the

284 marine diatom Phaeodactylum tricornutum. *Genome Biol* 2015;16:102.

285 Zang, C*., et al.* A clustering approach for identification of enriched domains from histone

286    modification ChIP-Seq data. *Bioinformatics* 2009;25(15):1952-1958.

287    Zhang, Y.*, et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol* 2008;9(9):R137.

288    Zhao, X.*, et al.* Genome wide natural variation of H3K27me3 selectively marks genes predicted

289    to be important for cell differentiation in Phaeodactylum tricornutum. *New Phytol*

290    2021;229(6):3208-3220.

291

292

293


294    **Legend**

295

296    **Figure 1**. Snapshot of PhaeoEpiView browser illustrating the different tracks of genes,
297    transposable elements, histone marks and DNA methylation. Both peaks and log2 fold
298    enrichment between IP and Input are displayed for H3K27me3.

299

300    **Supplementary Figure 1**. Comparison of RNA-seq quantification per gene on the 2008 (33
301    scaffold/chromosomes) and 2021 (25 chromosomes) assembled genomes. Two RNA-seq
302    replicates were used.

303

304    **Supplementary Table 1**. List of the genes not recovered on the lifted annotation with their
305    coordinates and features.

306

307    **Supplementary File 1.** Supplementary materials and methods.

308

309    **Supplementary File 2.** GFF3 file annotating *Phaeodactylum tricornutum* 2021 assembly
310    with Phatr3 lifted genes.

311